

Disclosed is a method for character spacing in text recognition tasks, wherein the possible intersecting points of extraction objects being examined are detected by means of white space and angular analysis. Plausible spacing lines are determined on the basis of the intersecting points and the corresponding counterpoints. The objects thus spaced undergo a classification process, final spacing being determined by the results thereof.

(57) Zusammenfassung

Es wird ein Verfahren zur Zeichentrennung bei Texterkennungsaufgaben angegeben, bei dem zu den untersuchten Extraktionsobjekten mittels Weißdellenanalyse und Winkelanalyse mögliche Schnittpunkte ermittelt werden, daß aus den Schnittpunkten und entsprechenden Gegenpunkten plausible Trennlinien ermittelt werden und daß die solcherart getrennten Objekte Klassifikationsverfahren unterzogen werden und auf der Grundlage der Ergebnisse die endgültige Trennung erfolgt.

LEDIGLICH ZUR INFORMATION

Codes zur Identifizierung von PCT-Vertragsstaaten auf den Kopfbögen der Schriften, die internationale Anmeldungen gemäss dem PCT veröffentlichen.

AL	Albanien	ES	Spanien	LS	Lesotho	SI	Slowenien
AM	Armenien	FI	Finnland	LT	Litauen	SK	Slowakei
AT	Österreich	FR	Frankreich	LU	Luxemburg	SN	Senegal
AU	Australien	GA	Gabun	LV	Lettland	SZ	Swasiland
AZ	Aserbaidsschan	GB	Vereinigtes Königreich	MC	Monaco	TD	Tschad
BA	Bosnien-Herzegowina	GE	Georgien	MD	Republik Moldau	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagaskar	TJ	Tadschikistan
BE	Belgien	GN	Guinea	MK	Die ehemalige jugoslawische Republik Mazedonien	TM	Turkmenistan
BF	Burkina Faso	GR	Griechenland	ML	Mali	TR	Türkei
BG	Bulgarien	HU	Ungarn	MN	Mongolei	TT	Trinidad und Tobago
BJ	Benin	IE	Irland	MR	Mauretanien	UA	Ukraine
BR	Brasilien	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Island	MX	Mexiko	US	Vereinigte Staaten von Amerika
CA	Kanada	IT	Italien	NE	Niger	UZ	Usbekistan
CF	Zentralafrikanische Republik	JP	Japan	NL	Niederlande	VN	Vietnam
CG	Kongo	KE	Kenia	NO	Norwegen	YU	Jugoslawien
CH	Schweiz	KG	Kirgisistan	NZ	Neuseeland	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Demokratische Volksrepublik Korea	PL	Polen		
CM	Kamerun	KR	Republik Korea	PT	Portugal		
CN	China	KZ	Kasachstan	RO	Rumänien		
CU	Kuba	LC	St. Lucia	RU	Russische Föderation		
CZ	Tschechische Republik	LJ	Liechtenstein	SD	Sudan		
DE	Deutschland	LK	Sri Lanka	SE	Schweden		
DK	Dänemark	LR	Liberia	SG	Singapur		
EE	Estland						

Verfahren zur Zeichentrennung bei Texterkennungsaufgaben

Die Erfindung betrifft ein Verfahren zur Zeichentrennung bei Texterkennungsaufgaben.

5

Bei der automatischen Erkennung von Texten, d.h. bei der Umwandlung der grafischen Information eines Dokumentes in Textzeichen, die mittels elektronischen Textverarbeitungsprogrammen weiterverarbeitet werden können, ist eine wesentliche Voraussetzung für einen erfolgreichen Erkennungsvorgang die genaue Bestimmung der Lage und der Größe der einzelnen Zeichen. Diese Bestimmung ist bei Vorlagen mit schlechtem Schriftbild oder Schriftarten mit sehr engem Zeichenabstand unter anderem dadurch problematisch, daß die Zeichen mit einander verbunden sind „zusammenwachsen“ und damit durch herkömmliche Methoden wie die einfache Konturverfolgung nicht mehr getrennt werden können.

20 Der Erfindung liegt daher die Aufgabe zugrunde, ein verbessertes Verfahren zur Trennung miteinander verbundener Zeichen anzugeben.

Dies geschieht erfindungsgemäß mit einem Verfahren der eingangs genannten Art, bei dem zu den untersuchten Extraktionsobjekten mittels Weißdellenanalyse und Winkelanalyse mögliche Schnittpunkte ermittelt werden, bei dem aus den Schnittpunkten und entsprechenden Gegenpunkten plausible Trennlinien ermittelt werden und bei dem die solcherart getrennten Objekte Klassifikationsverfahren unterzogen werden und auf der Grundlage der Ergebnisse die endgültige Trennung erfolgt.

35 Vorteilhaft ist eine Ausgestaltung des Verfahrens in der Weise, daß bei mehr als drei möglichen Schnittpunkten, ein erster Schnitt durch den vom linken Zeichenanfang gezählten

vierten Schnittpunkt erfolgt. Dies deswegen, weil kein übliches Textzeichen der lateinischen Schrift mehr als drei Weißdellen aufweist.

5 Günstig ist es ferner, wenn nach einem ersten Schnitt mit einem ersten möglichen Schnittpunkt und einem darauffolgenden erfolglosen Klassifikationsversuch als Basis für einen weiteren Trennversuch der zum ersten möglichen Schnittpunkt nächstliegende linke Nachbarschnittpunkt vorgesehen wird.

10 Die Erfindung wird anhand von Figuren näher erläutert. Es zeigen beispielhaft:

Fig.1 eine Darstellung zur Weißdellenanalyse eines Bildes, Fig.2 eine Darstellung zur eigentlichen Zeichentrennung.

15 Der Ablauf des erfindungsgemäßen Verfahrens ist wie folgt:

Das Verfahren wird im Erkennungsvorgang nach der Bestimmung der Lage der Zeile gestartet. Bei der Ermittlung des Umfanges
20 eines Zeichens oder mehrerer verbundener Zeichen durch Konturverfolgung wird bereits eine Weißdellenanalyse durchgeführt. Nach dem Vorliegen der vollständigen Kontur erfolgt eine Winkelanalyse.

25 Mittels Weißdellenanalyse und Winkelanalyse werden mögliche Schnittpunkte ermittelt, die in Verbindung mit Gegenpunkten mögliche Trennlinien liefern.

Die Schnittpunkte werden hinsichtlich ihrer Plausibilität
30 untersucht. Dabei wird ermittelt, welche Zeichenfolgen die vorliegende Weißdellenkombination beinhalten. So sind beispielsweise in der Buchstabenfolge **WV** folgende Weißdellen enthalten OBEN-UNTEN-OBEN-UNTEN-OBEN. Wobei OBEN (UNTEN) eine nach oben (unten) offene Weißdelle kennzeichnet. Aus der
35 Kenntnis der Buchstaben heraus wird nun die erste Trennung durch den Schnittpunkt der vierten Weißdelle erfolgen.

Darauf wird ermittelt, inwieweit die Trennung des Objekts entlang der auf plausiblen Schnittpunkten beruhenden Trennlinien zu plausiblen Klassifikationsergebnissen führt. Mit anderen Worten, die getrennten Zeichen oder Zeichenteile werden einem Erkennungsvorgang z.B. mittels neuronalem Netz unterworfen und wenn dieser Vorgang zu einem zufriedenstellenden Ergebnis - einem mit hoher Sicherheit erkannten Zeichen - führt, dann wird die Trennung akzeptiert. Andernfalls wird die Trennung entlang von anderen Trennlinien solange wiederholt, bis ein zufriedenstellendes Ergebnis vorliegt.

Neuronale Netze sind mathematische Modelle, welche dem Aufbau des menschlichen Gehirns nachempfunden sind. Sie bestehen aus Neuronen, das sind im wesentlichen Summierelemente mit gewichteten Eingängen und einem nichtlinearen Verstärkeranteil, die zu einem parallelen Netzwerk mit typisch zwei Ebenen zusammengefaßt werden. Eine ausführliche Beschreibung des beim Ausführungsbeispiel eingesetzten „Feedforward Neural Networks“ findet sich beispielsweise in „Layered Neural Nets for Pattern Recognition“, B. Widrow, R. G. Winter, R. A. Baxter; IEEE Transactions on Acoustics, Speech and Signal Processing, Vol. 36. No. 7. July 88.

Die Mustererkennung mittels neuronalem Netz erfolgt nach dem in „A rotation, scaling, and translation invariant pattern classification system“, C. Yüceer, K Oflazer; Pattern Recognition, Vol. 26, No5, pp687-710, 1993. beschriebenen Verfahren.

30

Die Weißdellenanalyse wird anhand der Figur 1 näher beschrieben. Die Figur zeigt die beiden miteinander verbundenen Buchstaben x und f die eine Weißdelle W aufweisen. Unter Weißdelle W wird dabei ein von drei Seiten begrenzter weißer Zwischenraum verstanden, der eine gewisse

35

Tiefe aufweist und dessen offene Seite nach oben oder unten gerichtet ist. Ermittelt wird diese Weißdelle W bei der Verfolgung der Kontur des (zusammengewachsenen) Zeichens wenn die Konturlinie C zwei vorgegebene Schwellwerte SW in beiden
5 Richtungen überschreitet. Liegt wie in dem Beispiel eine nach unten offene Weißdelle W vor, dann wird der höchste Punkt der Konturlinie C als möglicher Schnittpunkt S definiert, bei einer nach oben offenen Weißdelle ist dies der tiefste Punkt.

10 Der Ablauf der daraufhin erfolgenden Winkelanalyse ist wie folgt:

Aus jeweils drei Punkten der Konturlinie C[i] werden zwei Vektoren ermittelt, für die gilt:

$$\vec{A} = C[i]C[i-5] \text{ und } \vec{B} = C[i]C[i+5]$$

15

Der Winkel zwischen den beiden Vektoren wird berechnet. Ist dieser linksläufig, mit einem Betrag kleiner als 80° und einer entweder nach oben oder nach unten weisenden Spitze (C[i]), dann wird der Winkel in eine Liste eingetragen.

20

Ist diese Bedingung für mehrere nebeneinander liegende Vektorpaare erfüllt, dann wird nur der Winkel mit dem geringsten Betrag weiterverfolgt.

25 Die in der Liste eingetragenen Winkel werden nun daraufhin untersucht, ob auf der gegenüberliegenden Seite der Konturlinie ein Winkel mit entgegengesetzter Orientierung der Spitze vorhanden ist. Ist dies der Fall, dann wird das daraufhin gebildete Winkelpaar als Position eines möglichen
30 Schnittpunktes gespeichert.

Im Folgenden der Ablauf bei der Bestimmung des Winkels zwischen zwei Vektoren, die durch 3 Punkte aus der Konturlinie($C_1:x_1/y_1$, $C_6:x_2/y_2$, $C_{11}:mx/my$) definiert. Daraus werden die x und y Komponenten der beiden Vektoren ermittelt.

$$Ax = x_1 - mx ; Ay = y_1 - my ; Bx = x_2 - mx ; By = y_2 - my ;$$

Der Winkel zwischen den Vektoren A und B wird wie folgt berechnet: Zuerst wird der Winkel von A zur x-Achse und dann der Winkel B zur x-Achse ermittelt.

$$Winkel = \arccos\left(\frac{\bar{A}x}{\sqrt{(\bar{A}x)^2 + (\bar{A}y)^2}}\right)$$

$$Winkel(inGrad) = \frac{Winkel(inRad) * 180}{Pi}$$

Winkel = 360 - winkelB + winkelA (ist Winkel größer 360°, dann wird der Winkel um 360° korrigiert)

Die Bestimmung der Winkelspitzenrichtung beruht auf der Überlegung, daß bei einer nach unten gerichteten Spitze die Y-Koordinaten der Punkte C_1 und C_6 kleiner als die Y-Koordinate von C_{11} sind.

Bei einer nach oben gerichteten Spitze müssen hingegen die Y-Koordinaten der Punkte C_1 und C_6 größer als die Y-Koordinate von C_{11} sein.

Die Eigenheiten gedruckter Texte und der Einfluß der begrenzten Bildauflösung bringen es mit sich, daß im Bereich eines Knicks der Kontur eines Zeichens die in der beschriebenen Weise ermittelten Winkel zwischen 2 Vektoren

abhängig vom Betrachtungsraum zuerst zunehmend kleiner werden und danach wieder kontinuierlich zunehmen. Für die weitere Auswertung wird daher nur der jeweils minimale Winkel eines derartigen Bereiches verwendet.

5

Zur Festlegung einer möglichen Trennlinie muß nun zu jedem möglichen Schnittpunkt $C(Nr)$ ein entsprechender Gegenpunkt auf dem gegenüberliegenden Zweig der Konturlinie $C(i); i=(0, \dots, \text{contourNr})$ ermittelt werden.

10

Dazu wird eine Gerade durch zwei auf der Konturlinie dem möglichen Schnittpunkt $C(Nr)$ benachbarte Punkte $C(Nr-1)$ und $C(Nr+1)$ gelegt, und zu dieser Geraden die Normale ermittelt. Die zu dem Schnittpunkt dieser Normalen mit dem

15 gegenüberliegenden Zweig der Konturlinie benachbarten Punkte werden hinsichtlich ihres Abstandswertes zum möglichen Schnittpunkt und der Normalen untersucht und der Konturpunkt mit dem minimalen Abstandswert als Gegenpunkt $C(g)$ und damit als zweiter Punkt der möglichen Trennlinie definiert.

20 Die mathematische Definition dieses Vorganges lautet:

$$nx = C(Nr+1)x - C(Nr-1)x$$

$$ny = C(Nr+1)y - C(Nr-1)y$$

$$\text{Abstand} = \sqrt{(C(Nr)x - C(i)x)^2 + (C(Nr)y - C(i)y)^2}$$

$$= \text{abs} \left(\frac{nx * (C(i)x - C(Nr)x) + ny * (C(i)y - C(Nr)y)}{\sqrt{nx^2 + ny^2}} \right)$$

25 Abstand zu $g2$

$$\text{Abstandswert} = \text{Abstand} + \text{Abstand zu } g2;$$

$$C(g) = C(i) \mid \text{Abstandswert}(C(g), C(Nr)) = \min$$

Die eigentliche Trennung wird anhand der Fig. 2 erläutert:
Basis der Trennung ist die Konturlinie der extrahierten
Zeichen. In einem 1. Schritt wird ein Trennlinienpuffer mit 0
initialisiert, dies entspricht einer senkrechten Linie am
5 linken Rand, danach wird der am weitesten rechts liegende
Punkt der Konturlinie 1 zwischen 0 und dem der Trennung
zugrundeliegenden Schnittpunkt (das X-Wert-Maximum)
ermittelt. Ebenso werden der am weitesten rechts liegende
Punkt des Zweiges (das x-Wert Maximum) der Konturlinie vom
10 Gegenpunkt bis zum Ende der Kontur 2 und der Trennlinie 3
ermittelt.

Die gesammelten maximalen x-Werte stellen also den äußersten
rechten Rand des zur Klassifizierung herangezogenen Zeichens
15 dar.

Patentansprüche

- 1) Verfahren zur Zeichentrennung bei Texterkennungsaufgaben,
dadurch gekennzeichnet, daß zu den untersuchten
5 Extraktionsobjekten mittels Weißdellenanalyse und
Winkelanalyse mögliche Schnittpunkte ermittelt werden, daß
aus den Schnittpunkten und entsprechenden Gegenpunkten
plausible Trennlinien ermittelt werden und daß die solcherart
getrennten Objekte Klassifikationsverfahren unterzogen werden
10 und auf der Grundlage der Ergebnisse die endgültige Trennung
erfolgt.
- 2) Verfahren nach Anspruch 1, **dadurch gekennzeichnet**, daß bei
mehr als drei möglichen Schnittpunkten , ein erster Schnitt
15 durch den vom linken Zeichenanfang gezählten vierten
Schnittpunkt erfolgt.
- 3) Verfahren nach Anspruch 1 oder 2, **dadurch gekennzeichnet**,
daß nach einem ersten Schnitt mit einem ersten möglichen
20 Schnittpunkt und einem darauffolgenden erfolglosen
Klassifikationsversuch als Basis für einen weiteren
Trennversuch der zum ersten möglichen Schnittpunkt
nächstliegende linke Nachbarschnittpunkt vorgesehen wird.

1/1

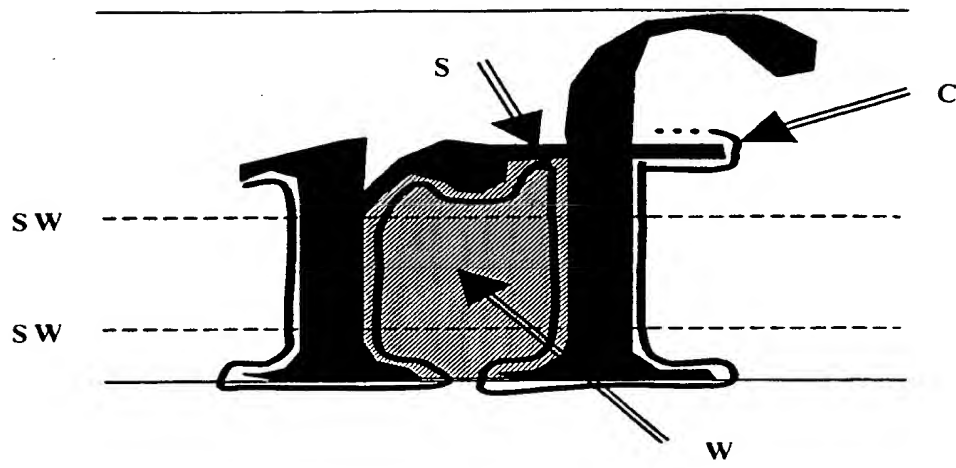


Fig. 1

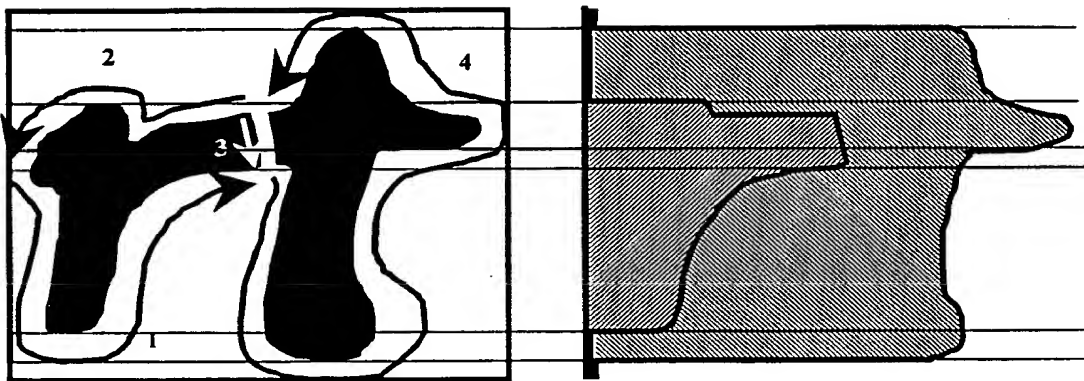


Fig. 2

THIS PAGE BLANK (USPTO)

INTERNATIONAL SEARCH REPORT

International Application No

PCT/EP 99/06841

A. CLASSIFICATION OF SUBJECT MATTER
IPC 7 G06K9/34

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 7 G06K

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	ZHERU CHI ET AL: "SEPARATION OF SINGLE- AND DOUBLE-TOUCHING HANDWRITTEN NUMERAL STRINGS" OPTICAL ENGINEERING, US, SOC. OF PHOTO-OPTICAL INSTRUMENTATION ENGINEERS. BELLINGHAM, vol. 34, no. 4, page 1159-1165 XP000497484 ISSN: 0091-3286 the whole document	1

☐ Further documents are listed in the continuation of box C.

☐ Patent family members are listed in annex.

* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"S" document member of the same patent family

Date of the actual completion of the international search

20 December 1999

Date of mailing of the international search report

11/01/2000

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax (+31-70) 340-3016

Authorized officer

Granger, B

THIS PAGE BLANK (USPTO)

INTERNATIONALER RECHERCHENBERICHT

Internationales Aktenzeichen

PCT/EP 99/06841

A. KLASSIFIZIERUNG DES ANMELDUNGSGEGENSTANDES

IPK 7 G06K9/34

Nach der internationalen Patentklassifikation (IPK) oder nach der nationalen Klassifikation und der IPK

B. RECHERCHIERTE GEBIETE

Recherchierte Mindestprüfstoff (Klassifikationssystem und Klassifikationssymbole)

IPK 7 G06K

Recherchierte aber nicht zum Mindestprüfstoff gehörende Veröffentlichungen, soweit diese unter die recherchierten Gebiete fallen

Während der internationalen Recherche konsultierte elektronische Datenbank (Name der Datenbank und evtl. verwendete Suchbegriffe)

C. ALS WESENTLICH ANGESEHENE UNTERLAGEN

Kategorie*	Bezeichnung der Veröffentlichung, soweit erforderlich unter Angabe der in Betracht kommenden Teile	Betr. Anspruch Nr.
X	ZHERU CHI ET AL: "SEPARATION OF SINGLE- AND DOUBLE-TOUCHING HANDWRITTEN NUMERAL STRINGS" OPTICAL ENGINEERING, US, SOC. OF PHOTO-OPTICAL INSTRUMENTATION ENGINEERS. BELLINGHAM, Bd. 34, Nr. 4, Seite 1159-1165 XP000497484 ISSN: 0091-3286 das ganze Dokument	1

☐ Weitere Veröffentlichungen sind der Fortsetzung von Feld C zu entnehmen

☐ Siehe Anhang Patentfamilie

* Besondere Kategorien von angegebenen Veröffentlichungen :

"A" Veröffentlichung, die den allgemeinen Stand der Technik definiert, aber nicht als besonders bedeutsam anzusehen ist

"E" älteres Dokument, das jedoch erst am oder nach dem internationalen Anmeldedatum veröffentlicht worden ist

"L" Veröffentlichung, die geeignet ist, einen Prioritätsanspruch zweifelhaft erscheinen zu lassen, oder durch die das Veröffentlichungsdatum einer anderen im Recherchenbericht genannten Veröffentlichung belegt werden soll oder die aus einem anderen besonderen Grund angegeben ist (wie ausgeführt)

"O" Veröffentlichung, die sich auf eine mündliche Offenbarung, eine Benutzung, eine Ausstellung oder andere Maßnahmen bezieht

"P" Veröffentlichung, die vor dem internationalen Anmeldedatum, aber nach dem beanspruchten Prioritätsdatum veröffentlicht worden ist

"T" Spätere Veröffentlichung, die nach dem internationalen Anmeldedatum oder dem Prioritätsdatum veröffentlicht worden ist und mit der Anmeldung nicht kollidiert, sondern nur zum Verständnis des der Erfindung zugrundeliegenden Prinzips oder der ihr zugrundeliegenden Theorie angegeben ist

"X" Veröffentlichung von besonderer Bedeutung; die beanspruchte Erfindung kann allein aufgrund dieser Veröffentlichung nicht als neu oder auf erfinderischer Tätigkeit beruhend betrachtet werden

"Y" Veröffentlichung von besonderer Bedeutung; die beanspruchte Erfindung kann nicht als auf erfinderischer Tätigkeit beruhend betrachtet werden, wenn die Veröffentlichung mit einer oder mehreren anderen Veröffentlichungen dieser Kategorie in Verbindung gebracht wird und diese Verbindung für einen Fachmann naheliegend ist

"Z" Veröffentlichung, die Mitglied derselben Patentfamilie ist

Datum des Abschlusses der internationalen Recherche

20. Dezember 1999

Abschließdatum des internationalen Recherchenberichts

11/01/2000

Name und Postanschrift der internationalen Recherchenbehörde

Europäisches Patentamt, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax (+31-70) 340-3018

Bevollmächtigter Beauftragter

Granger, B

THIS PAGE BLANK (USPTO)